



# Quantifier elimination in automatic loop parallelization

Armin Größlinger, Martin Griebel, Christian Lengauer\*

*Department of Mathematics and Informatics, University of Passau, Germany*

Received 9 December 2004; accepted 23 September 2005

Available online 7 September 2006

---

## Abstract

We present an application of quantifier elimination techniques in the automatic parallelization of nested loop programs. The technical goal is to simplify affine inequalities whose coefficients may be unevaluated symbolic constants. The values of these so-called structure parameters are determined at run time and reflect the problem size. Our purpose here is to make the research community of quantifier elimination, in a tutorial style, aware of our application domain – loop parallelization – and to highlight the role of quantifier elimination, as opposed to alternative techniques, in this domain. Technically, we focus on the elimination method of Weispfenning.

© 2006 Elsevier Ltd. All rights reserved.

**Keywords:** Loop parallelization; Quantifier elimination

---

## 1. Introduction

Loop parallelization is a problem which occurs predominantly in high-performance computing. It is applied to highly iterative, computationally-intensive matrix algorithms in the fields of linear algebra, numerical computation, image, signal and text processing, computational physics, chemistry and biology, etc. The algorithms are written as sequential programs and parallelism is a device to speed up the computation dramatically, possibly using hundreds or thousands of processors.

The *polytope model* (Lengauer, 1993; Feautrier, 1996) is a powerful geometric model for the parallelization of nested loop programs. Since its inception, it has been extended from polytopes to polyhedra, and is often also called the *polyhedron model*.

---

\* Corresponding address: University of Passau, Room 129, Innstrasse 33, 94030 Passau, Germany.

E-mail addresses: [groessli@fmi.uni-passau.de](mailto:groessli@fmi.uni-passau.de) (A. Größlinger), [griebel@fmi.uni-passau.de](mailto:griebel@fmi.uni-passau.de) (M. Griebel), [lengauer@fmi.uni-passau.de](mailto:lengauer@fmi.uni-passau.de) (C. Lengauer).

```

DO i=0,n
  DO j=0,n
    C(i+j) = C(i+j) + A(i) * B(j)
  
```

Fig. 1. A sequential source code for the polynomial product.

The parallelization takes place in three steps:

- 1. Modelling:** The sequential input program, typically but not necessarily written in Fortran or C, is transformed to a model-based description.
- 2. Parallelization:** The parallelization converts the source model to a target model.
- 3. Code generation:** The code generation transforms the target model to an executable (parallel) program.

The following sections describe these steps in more detail and show how traditional methods of the polyhedron model solve the different mathematical problems which arise — provided they can be described with linear formulas. Then, we illustrate, in a tutorial style, how the application of the quantifier elimination method by Weispfenning (Weispfenning, 1988; Loos and Weispfenning, 1993; Weispfenning, 1994b, 1997) can help to overcome this proviso.

### 1.1. Modelling

The steps of a loop nest, iterating over an array structure, are laid out in a multi-dimensional integer coordinate system, with one dimension per loop. The set of points which result is called the *index space*. The borders of the index space are given by the bounds of the loops. Thus, the index space is described by a set of affine inequalities — a standard description of a polyhedron.

Directed edges between the points in the index space represent the dependences between the loop iterations. There is a vibrant research area concerned with the technology for generating these dependences automatically from the source program (Feautrier, 1991; Collard and Griebel, 1999; Amiranoff, 2004).

*Example.* As an example, let us look at the core part of the simple sequential algorithm for the polynomial product (Fig. 1). This code computes the coefficients of the product (in array  $C$ ) from the coefficients of the two factors (in arrays  $A$  and  $B$ ). Both factors have  $n + 1$  coefficients;  $n$  is called a *structure parameter*. Dependences arise from the accesses to elements of the array  $C$ . Several iterations of the loop nest access the same memory location, since the expression  $i + j$  has the same value for different instances of  $i$  and  $j$ . Thus, all iterations  $(i_0, j_0)$  and  $(i_1, j_1)$  which satisfy the memory equation  $i_1 + j_1 = i_0 + j_0$  are in mutual dependence.

In order to keep the representation simple, we disregard all dependences that are a transitive combination of other dependences and consider only the rest: the so-called *spanning dependences*. Starting from an iteration  $(i_0, j_0)$ , we compute those iterations  $(i_1, j_1)$  that, first, precede  $(i_0, j_0)$ , second, satisfy the memory equation  $i_1 + j_1 = i_0 + j_0$ , and third, are maximal among all solutions  $(i_1, j_1)$  w.r.t. the lexicographic order. Solving this system of constraints, we obtain  $(i_1, j_1) = (i_0 - 1, j_0 + 1)$ .

Fig. 2 depicts the index space and the spanning dependences between loop iterations (denoted by arrows from the earlier to the later iteration accessing the same array element); the full set of dependences is the transitive closure of the dependences shown.

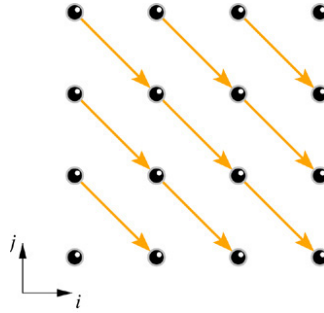


Fig. 2. Dependence structure of the polynomial product.

*Limitations.* The polyhedron model comes with restrictions:

- The program must be a loop nest whose body consists of statements. The loop nest may be *imperfect*, i.e., not all statements need be in the body of the innermost loop.
- Statements are allowed to be array assignments, conditionals and subprogram calls. In the parallelization, the latter are treated as read and write accesses to all parameters.
- The loop bounds and array index expressions must be affine and may contain other variables which are constant in the loop nest.

### 1.2. Parallelization

With techniques of linear algebra and linear programming, an automatic, optimizing search yields the best mapping of the loop steps to time and space (processors) with respect to some objective function like the number of parallel execution steps (the most popular choice), the number of communications, a balanced processor load or combinations of these or others. When applied to the source index space, this so-called *space–time mapping* specifies another polyhedral space, the so-called *target index space*.

Methods based on the polyhedron model have been implemented in various prototypical preprocessors. One for C with the message-passing library MPI is our own loop restructurer LooPo (Griebl and Lengauer, 1996). These systems use a number of well known *schedulers* (which compute temporal distributions) (Feautrier, 1992a,b) and *allocators* (which compute spatial distributions) (Feautrier, 1994; Dion and Robert, 1995) for the optimized search of a space–time mapping. In other words, for each computation of each statement in the loop nest, the scheduler determines its logical time step and the allocator its logical processor. The scheduler must ensure that the execution date increases in the direction of each dependence; this constraint can be formulated as a finite system of inequalities that the scheduler has to solve. The allocator tries to optimize locality by putting computations, which access the same memory cell, onto the same logical processor; this is done by a heuristic search.

*Example.* For the polynomial product (Figs. 1 and 2), the top-most row of iterations in Fig. 2 ( $j = n$ ) can be executed first and simultaneously. Once it is computed, all iterations of the second row can be computed, and so on. Thus, a possible time mapping is  $t(i, j) = n - j$ . This result is also found by solving the scheduler's constraint system. (Another solution is  $t'(i, j) = i$ .)

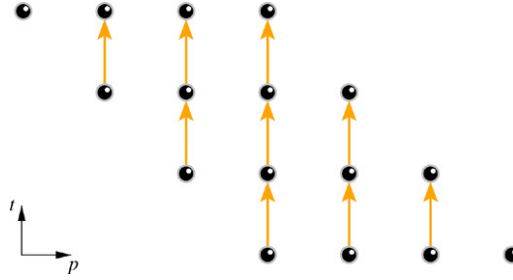


Fig. 3. Transformed index space and dependences of the polynomial product.

Since the allocator tries to put dependent computations on the same processor, one possible choice of time mapping  $t(i, j)$  and space mapping  $p(i, j)$  is:

$$\begin{aligned} t(i, j) &= n - j \\ p(i, j) &= i + j. \end{aligned} \tag{1}$$

The target space and the dependences in the target space derived from these mappings are depicted in Fig. 3.

*Limitations.* The only restriction w.r.t. the parallelization phase is that the space–time mapping must be affine.

Recent extensions allow mild violations of this affinity requirement. Essentially, they permit a constant number of breaks in the affinity (e.g., a time mapping like  $t''(i, j) = \text{if } i < n - j \text{ then } i \text{ else } n - j$  for our running example).

### 1.3. Code generation

The biggest challenge is to convert the solutions found in the model into efficient code. Significant headway has been made recently on how to avoid frequent run-time tests which guide control through the various parts of the iteration space (Quilleré et al., 2000; Bastoul, 2004).

The result of the code generation is a loop nest whose loops are sequential or parallel depending on whether the corresponding axis in the target coordinate system represents time or space.

While the axes of the target coordinate system are unordered, the loops in the target loop nest must be ordered. Different orders can be chosen. If the outer loop is sequential, one obtains *synchronous* parallelism, i.e., all processes iterate following the same global clock. If the outer loop is parallel, one obtains *asynchronous* parallelism, i.e., each process iterates following its own, local clock. However, the bound expressions of every loop may only contain iteration variables of the outer loops. To enforce this in the code generation, one must use a process of variable elimination, which is in essence a quantifier elimination problem. Care must be taken that the system resulting from the variable elimination has a syntactic form suitable for code generation (see Section 2 for more detail). In simple cases, e.g., for loop nests with only one statement, Fourier–Motzkin elimination can be used, but the general case requires more advanced techniques (Quilleré et al., 2000; Bastoul, 2004).

*Example.* From the space–time mapping, an asynchronous target program can be synthesized, which is shown in Fig. 4. To generate the asynchronous target program, the target index space

```

DO PAR p=0,2*n
  DO t=max(0,n-p),min(n,2*n-p)
    C(p) = C(p) + A(t+p-n) * B(n-t)

```

Fig. 4. A parallel target code for the polynomial product.

(Fig. 3) described by the inequalities:

$$\begin{aligned}
 0 &\leq t \leq n \\
 n - t &\leq p \leq 2 \cdot n - t
 \end{aligned}$$

has to be expressed differently. To generate code, i.e., for-loops which enumerate every integral point in the target index space, where  $p$  is the variable of the outer loop and  $t$  the variable of the inner loop, we have to express the target space as follows:

$$\begin{aligned}
 0 &\leq p \leq 2 \cdot n \\
 0 &\leq t \leq n \\
 n - p &\leq t \leq 2 \cdot n - p
 \end{aligned} \tag{2}$$

to derive the bounds of the for-loops from the system (compare system (2) to the loop bounds in Fig. 4).

#### 1.4. An extension: Granularity control

Methods based on the polyhedron model are elegant and work well. However, as presented so far, the resulting parallelism is usually too fine-grained to be efficient. An extension to the model-based parallelization phase, called *tiling*, solves this granularity problem (Andonov et al., 2001; Griehl et al., 2004). It partitions the index space by covering it with equally shaped and sized polyhedra, mostly parallelepipeds. Every coordinate of a point in the index space is divided into two components: the coordinate of the tile, and the coordinate of the local offset within the tile.

There is one limitation of the descriptive power of the current polyhedron model, which is hindering the applicability of the tiling technique. If the target program is not to be compiled for a fixed number of processors (which also fixes the size of the tiles), but for an unknown, i.e., parametric number of processors, the size of the tiles depends on parameters. Such a tiling can only be described by an inequality system in which the respective parameters appear as coefficients of variables. Unfortunately, the tools available for solving the linear algebra and linear optimization problems do not handle this case; they require coefficients which are numeric constants.

We have developed techniques for handling parametric coefficients in the polyhedron model. Our application is parametric tiling but, in principle, these techniques can also be used to generalize algorithms in other areas of the polyhedron model, e.g., the dependence analysis.

*Example.* The target code for the polynomial product shown in Fig. 4 has two drawbacks: it requires  $2 \cdot n + 1$  processors and the work is unevenly distributed among them. With tiling, we can solve both problems at the same time, i.e., balance the work among the processors by combining smaller into bigger chunks of work, and map several of the  $2 \cdot n + 1$  “virtual” processors onto one physical processor.

In order to achieve the desired, so-called *block-cyclic* tiling, we express the virtual processor coordinate  $p$  by a coordinate triple  $(r, b, o)$ ; Fig. 6 depicts the situation for two real processors.

```

DO PAR r=0,NP-1
DO t=0,n
DO s=ceiling(-t/(NP*B) -r/NP +(n-B+1)/(NP*B)),
    floor(-t/(NP*B) -r/NP +(2*n)/(NP*B))
DO o=max(0,-t-r*B-s*NP*B+n),min(B-1,-t-r*B-s*NP*B+2*n)
DO p=max(n-t,r*B+s*NP*B+o),min(2*n-t,r*B+s*NP*B+o)
    C(p) = C(p) + A(t+p-n) * B(n-t)

```

Fig. 5. Target code for the polynomial product after tiling.

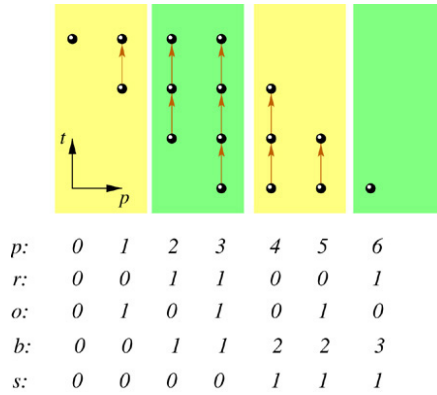


Fig. 6. Transformed index space after a block-cyclic tiling.

A number of adjacent index points are grouped together into a so-called block (in Fig. 6, the number is 2). The block's number is the coordinate  $b$ . Within each block, an index point has a unique offset, denoted by  $o$ . The blocks are distributed across the real processors in cyclic, i.e., alternating fashion. The assigned number of the real processor is coordinate  $r$ .

A block-cyclic tiling of a target space can always be imposed by adding the following set of inequalities to the description of the target space:

$$\begin{aligned}
 0 &\leq r \leq NP - 1 \\
 0 &\leq o \leq B - 1 \\
 p &= b \cdot B + o \\
 b &= s \cdot NP + r
 \end{aligned} \tag{3}$$

where  $NP$  is the number of real processors available and  $B$  the block size of the block-cyclic tiling. The last equation (with variable  $s$ ) specifies the cyclic distribution of blocks across real processors, where  $s$  ranges over the cycles.

In order to generate a loop nest, we subject this (non-linear) system to a version of Fourier-Motzkin elimination that can deal with non-linear parameters (cf. applications in Section 2.2.1). The resulting loop code is given in Fig. 5.

As desired, this program ensures that a fixed number of processors is used (in Fig. 6:  $NP=2$ , counted by  $r$  and indicated by the shadings), and that the load is balanced (in Fig. 6: each of the two physical processors performs two computations per time step).

## 2. Quantifier elimination in the polyhedron model

Quantifier elimination is used in several places in the polyhedron model. Let us first look at its application in the current polyhedron model, and at the algorithms used (Section 2.1). Then, let us look at our generalized model, which allows non-linear parameters, and how the need for non-linear quantifier elimination arises.

We can distinguish the following applications on linear inequality systems:

- testing the feasibility,
- projecting,
- computing solutions.

*Feasibility testing.* Testing the feasibility of an inequality system means deciding the formula  $\varphi = \exists x_1 \dots \exists x_n(\psi)$ , where  $\psi$  describes the inequality system, i.e., is a conjunction of inequalities. Sometimes, it suffices to decide whether  $\mathbb{R} \models \varphi$  holds but, since the variables we analyze in the polyhedron model are integer-valued loop indices, knowing that a solution exists in the reals is often not sufficient. For example, dependence tests can output false dependences when they do not include a check that all the loop indices are integral.

*Projecting.* Projections of inequality systems are mainly needed during code generation. To generate code which enumerates every point in a polyhedron described by  $\psi$ , one performs quantifier elimination on the following formulas  $\varphi_1, \dots, \varphi_{n-1}$ :

$$\begin{aligned}\varphi_1 &= \exists x_2 \exists x_3 \dots \exists x_n(\psi) \\ \varphi_2 &= \exists x_3 \dots \exists x_n(\psi) \\ &\vdots \\ \varphi_{n-1} &= \exists x_n(\psi).\end{aligned}$$

This yields the quantifier-free equivalents  $\varphi'_1, \dots, \varphi'_{n-1}$ . Here,  $\varphi'_1$  is equivalent to a conjunction of atomic formulas describing the lower and upper bounds for the variable  $x_1$ ,  $\varphi'_2$  describes the lower and upper bounds for  $x_2$  in dependence of  $x_1$  (the conjuncts only involving  $x_1$  are ignored), and so on. Finally, in the formula  $\psi$ , the lower and upper bounds for  $x_n$  depend on  $x_1, \dots, x_{n-1}$ . The problem with using a quantifier elimination tool to compute  $\varphi'_1, \dots, \varphi'_{n-1}$  is that such tools usually do not guarantee that the results are conjunctions of atomic formulas, i.e., a more or less difficult post-processing effort (e.g., formula simplification) or the reliance on implementation details become necessary.

*Computing solutions.* Projecting an inequality system onto some of its dimensions can be viewed as computing the complete solution set of the inequality system. Sometimes, however, it is sufficient to find a single solution of a system, maybe under the constraint that the solution has some properties. For example, in dependence analysis, one can, instead of dealing with all dependences of some equivalence class, choose the lexicographic maximum as a representative (Feautrier, 1991).

### 2.1. The current polyhedron model

The current polyhedron model is restricted to inequality systems which can be described by formulas which are linear in the variables *and* in the structure parameters. The reason for this

restriction is the fact that subproblems encountered in many algorithms are again linear in the variables and in the structure parameters, so this ensures a kind of “closure” property.

*Feasibility testing.* Various algorithms and implementations thereof exist for testing the feasibility of an inequality in the reals. Most basically, Fourier–Motzkin elimination can decide the existence of a real solution. Because of its doubly exponential complexity, Fourier–Motzkin is usually replaced by other, simply exponential algorithms — most prominently, the Simplex algorithm, Chernikova’s algorithm (Chernikova, 1968), or Weispfenning’s algorithm for linear inequality systems (Weispfenning, 1994a). Deciding the existence of integral solutions is algorithmically harder and, therefore, fewer implementations exist. For conjunctions of inequalities, PIP (Parametric Integer Programming) can be used (Feautrier et al., 2003). For more general situations, the Omega test is available (Pugh and Wonnacott, 1992; Pugh, 1992), and there are computer algebra systems which handle integral quantifier elimination on linear formulas, e.g., the REDLOG package (Dolzmann and Sturm, 1997) in the computer algebra system REDUCE (Hearn, 2004).

*Projecting.* To compute the projection of an inequality system (the formulas  $\varphi'_1, \dots, \varphi'_{n-1}$ ), one does not use a full quantifier elimination tool, since there are algorithms which have been specifically designed to do so. For example, it is at the heart of Fourier–Motzkin to compute these formulas, or Chernikova’s algorithm can be used to compute the projections.

*Example.* For the polynomial product, the parallelization step (using the space–time mapping given in Eq. (1)) yields the target index space depicted in Fig. 3, which is described by the following inequality system:

$$\begin{aligned} 0 &\leq t \leq n \\ n - t &\leq p \leq 2 \cdot n - t. \end{aligned} \tag{4}$$

To obtain the (untiled) target code, we have to compute the inequalities shown in system (2). As described in Section 2, we can get the bounds for the inner loop index  $t$  directly from this inequality system (by solving each inequality for  $t$ ). The bounds for  $p$  can be obtained by computing a quantifier-free equivalent of the following predicate:

$$\exists t (0 \leq t \leq n \wedge n - t \leq p \leq 2 \cdot n - t).$$

Using Fourier–Motzkin elimination, we obtain the result  $0 \leq p \leq 2 \cdot n$  directly, and it is guaranteed that any result returned is a conjunction of atomic formulas. A full quantifier elimination procedure does not give such a guarantee. For example, REDLOG yields:

$$0 \leq p \leq n \quad \vee \quad n \leq p \leq 2 \cdot n.$$

Of course, this is equivalent to the result computed by Fourier–Motzkin, but it requires some suitable simplification procedure to be applied before code, i.e., a loop nest can be generated. For example, the formula simplifier SLFQ (Brown, 2002), which is based on QEPCAD, simplifies REDLOG’s result to  $0 \leq p \leq 2 \cdot n$ . Another way to find an equivalent conjunction is to check for every atomic formula  $\tau$  in the (negation free) result  $\varphi$  whether  $\mathbb{R} \models \varphi \rightarrow \tau$  holds, using quantifier elimination. A conjunction equivalent to  $\varphi$  is then given by the conjunction of all atomic formulas  $\tau$  for which  $\mathbb{R} \models \varphi \rightarrow \tau$  holds.



*Computing solutions.* To compute a single solution of an inequality system, we usually use the Simplex algorithm with a suitable target function; usually, the lexicographic minimum (as computed by PIP, for example) is the desired solution. Quantifier elimination with answer (which computes solutions for existentially quantifier variables) can also be used, but our experiments showed that, in the linear case, the specialized tools (like PIP) are faster, sometimes much faster. See Section 2.2.2 for an example application of quantifier elimination with answer.

## 2.2. The polyhedron model with non-linear parameters

As our example demonstrates, some desired transformations, e.g., tiling the target index space with a parametric tile size, cannot be described in the current polyhedron model, since they require parametric coefficients. This case is not handled by the algorithms mentioned in the last section. We have developed a generalization of the model to allow parameters in the coefficients of variables. The significant difference to the current model is that the most basic operation on equations and inequalities, namely solving for a variable, can, in general, not be performed uniformly for all parameter values, but a case distinction on the sign of the coefficient must be made, which can depend on the parameters. Therefore, algorithms which return a single answer in the current model (e.g., Fourier–Motzkin elimination) need to return a case distinction (with conditions on the parameters) and a separate solution for every case as result.

There are two ways of obtaining algorithms for the generalized model. First, we can look for a way to transform an existing algorithm into an algorithm for the generalized model; this is addressed in Section 2.2.1. Second, we can try to use a full quantifier elimination tool (which is not limited to formulas, which are linear in the variables and parameters, as are the tools mentioned in Section 2.1) to solve some problems; see Section 2.2.2. We review both approaches and refer the reader to our previous publications for the details (Größlinger et al., 2004; Größlinger, 2003).

Since the theory of integers with addition and multiplication is undecidable, we cannot, in general, apply the generalization technique presented in Section 2.2.1 to obtain algorithms for integral problems. Constructing such algorithms, which work for the cases usually encountered in the polyhedron model, is a subject of current research.

### 2.2.1. Generalizing algorithms by program transformation

The current polyhedron model uses a rich set of algorithms. Since some of them have been developed with the polyhedron model in mind, it is desirable to be able to transform these algorithms into algorithms which can handle non-linear parameters. In this section, we present our way to obtain a generalized algorithm by a transformation.

The algorithms used in the polyhedron model work on the coefficients of the inequalities; in the current model, they work on rational or integral numbers. As a very simple example of an algorithm which is used to manipulate the coefficients, let us consider the following function *abs* which computes the absolute value of its argument:

$$\begin{aligned} \text{abs} : \mathbb{Q} &\rightarrow \mathbb{Q} \\ \text{abs}(x) &= \begin{cases} x & \text{if } x \geq 0 \\ \text{abs}(-x) & \text{if } x < 0. \end{cases} \end{aligned}$$

We use recursion on purpose, since it turns out to be the problematic issue in the transformation.

For every argument  $x$  of *abs*, the conditions  $x \geq 0$  and  $x < 0$  can be evaluated and it can be determined whether to continue the computation with returning  $x$  or with computing  $\text{abs}(-x)$ .

In the generalized polyhedron model, the coefficients of the variables are polynomials or rational functions in the structure parameters. Therefore, in the generalized model, this algorithm has to deal with arguments  $f \in \mathbb{Q}[p_1, \dots, p_m]$ . But now,  $f \geq 0$  and  $f < 0$  cannot be evaluated, since their truth value depends on the parameters  $p_1, \dots, p_m$ , in general. The generalized version  $\text{abs}^g$  of  $\text{abs}$ , which we are about to construct, has to return a data structure which describes the value of  $\text{abs}(f)$  in dependence of the structure parameters. It seems that  $\text{abs}^g$  can be obtained very simply. Just replace the case distinction of  $x \geq 0$  vs.  $x < 0$  in the algorithm's *code* by a case distinction in the (generalized) algorithm's *result data structure*.

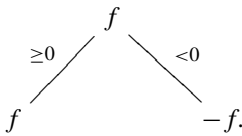
To represent case distinctions which return results of type  $\alpha$  as a data structure, we use decision trees, denoted by the data type  $\text{Tree } \alpha$ . This data type has (among others) the following constructors:

- $\text{Leaf} : \alpha \rightarrow \text{Tree } \alpha$  to represent a result of type  $\alpha$  as a decision tree.
- $\text{GeCond} : \mathbb{Q}[p_1, \dots, p_m] \times \text{Tree } \alpha \times \text{Tree } \alpha \rightarrow \text{Tree } \alpha$  to represent a binary case distinction.  $\text{GeCond}(g, t_1, t_2)$  means that, if  $g \geq 0$  (for some given values of the parameters), then  $t_1$  represents the result, otherwise  $t_2$ .

For example, we would like the absolute value of the polynomial  $f$  to be represented by the following expression:

$$\text{GeCond}(f, \text{Leaf } f, \text{Leaf } (-f)).$$

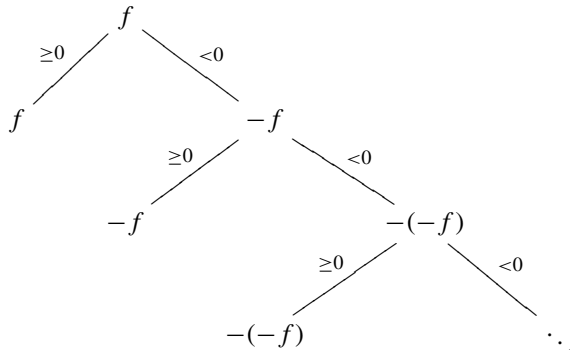
In our depiction of the decision tree, we label each edge with the condition which guards the subtree suspended from it (here, the subtrees are leaves):



Unfortunately, we must do some work to obtain this result. By transferring the case distinction from the code to the result data structure, we obtain the following implementation of  $\text{abs}^g$ :

$$\begin{aligned} \text{abs}^g &: \mathbb{Q}[p_1, \dots, p_m] \rightarrow \text{Tree } \mathbb{Q}[p_1, \dots, p_m] \\ \text{abs}^g(f) &= \text{GeCond}(f, \text{Leaf } f, \text{abs}^g(-f)). \end{aligned}$$

Alas, in the evaluation of  $\text{abs}^g(f)$ , the recursive call  $\text{abs}^g(-f)$  is evaluated repeatedly ad infinitum, which yields the following infinite tree:



We observe that, when  $f < 0$  holds,  $-f < 0$  cannot hold, but  $-f \geq 0$  must hold. i.e., the right child of the root node could be replaced by the leaf  $-f$ , but our transformed algorithm does not do so!

This example shows that a naïve transformation of algorithms to handle non-linear parameters is not sufficient, since recursive functions can produce infinite decision trees. Only after simplification (i.e., elimination of branches with contradictory conditions), the resulting decision tree may be finite. On the other hand, our transformational approach ensures that the results computed by the generalized algorithm are correct by construction, provided that the algorithm subject to the transformation is correct. Two questions must be answered:

- (1) How can the simplification of the decision trees be performed?
- (2) Does the simplification *guarantee* that only finite decision trees are constructed?

The answer to Question (1) is simple. Since the coefficients of the variables in the generalized polyhedron model are polynomials or rational functions in the parameters, and the algorithms manipulate the coefficients using addition, subtraction, multiplication, and division as algebraic operations and use the natural ordering relation, all the conditions in the decision trees can be written as  $f\rho 0$  where  $f \in \mathbb{Q}[p_1, \dots, p_m]$  and  $\rho \in \{=, \neq, <, \leq, \geq, >\}$ . Thus, real quantifier elimination can be used to decide the feasibility of the conjunction of the conditions guarding a branch (in the reals). We use REDLOG (Dolzmann and Sturm, 1997) (a quantifier elimination package for the computer algebra system REDUCE (Hearn, 2004)), which implements Weispfenning's elimination-by-virtual-substitution method (Weispfenning, 1988; Loos and Weispfenning, 1993; Weispfenning, 1997), and QEPCAD (Brown, 2004), which implements quantifier elimination through partial cylindrical algebraic decomposition. By performing a top-down search for contradictions this way, we avoid running into the infinite branch in the  $\text{abs}^g$  example, since the infinite branch is cut off when  $f < 0 \wedge -f < 0$  is found to be infeasible. Using real quantifier elimination, although the structure parameters are from  $\mathbb{Z}$ , turns out not to be a problem, since the structure of the problems (like tiling) implies that quantifier elimination in the reals is sufficient. For example, we have proved that tiling the index space of one statement with congruent parallelepipeds as tiles yields a target code which requires no case distinctions in the parameters at all (Gröbinger et al., 2004).

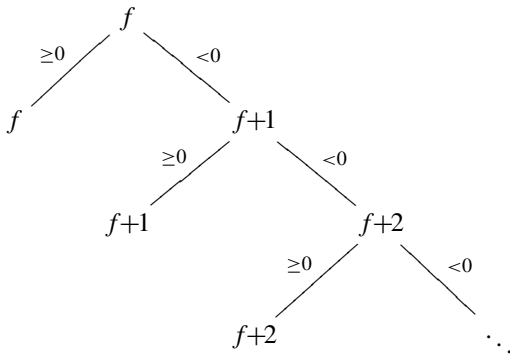
Question (2) is more difficult to answer. There are algorithms which produce infinite decision trees even with the best possible top-down tree simplification. For example, the following function:

$$\begin{aligned} \text{pos} : \mathbb{Q} &\rightarrow \mathbb{Q} \\ \text{pos}(x) &= \begin{cases} x & \text{if } x \geq 0 \\ \text{pos}(x + 1) & \text{if } x < 0 \end{cases} \end{aligned}$$

is generalized to:

$$\begin{aligned} \text{pos}^g : \mathbb{Q}[p_1, \dots, p_m] &\rightarrow \text{Tree } \mathbb{Q}[p_1, \dots, p_m] \\ \text{pos}^g(f) &= \text{GeCond}(f, \text{Leaf } f, \text{pos}^g(f + 1)) \end{aligned}$$

which produces (for some input  $f$ ) the following decision tree:



But every finite prefix of the infinite branch is consistent. So algorithm *pos* cannot be generalized by our transformation — even with simplification.

Termination (i.e., the construction of a finite decision tree) is guaranteed if every infinite branch has a finite prefix which is inconsistent. This condition ensures that searching top-down for contradictions and cutting off branches with infeasible conditions eliminates the infinite branch. Note that it does not suffice that, for any given values of the parameters, only a finite prefix of the tree is used. In *pos*<sup>g</sup>, this is the case — yet the tree cannot be made finite.

*Applications.* We have used this transformation to generalize Fourier–Motzkin elimination and the Simplex algorithm to non-linear parameters. For both algorithms, one can prove that top-down simplification ensures that only finite trees are generated. This gives us the ability to solve each of our three main applications (feasibility testing, projection, computing a solution).

Proving the termination of Fourier–Motzkin is rather simple. Since the number of case distinctions is bounded by the number of variables and inequalities in the input system, Fourier–Motzkin produces finite trees even without simplification. But simplification is still necessary to reduce the – in principle – exponential size of the decision trees to trees with few or no case distinctions in practical applications. E.g., for the tiling example, it turns out that no case distinction is ever needed.

The trees produced by the generalized Simplex are infinite, because they must also describe cyclic computations. But we can prove that the use of Bland’s minimal index rule for choosing pivots (Bland, 1977) ensures that top-down tree simplification cuts off every infinite branch. Assume that there is an infinite branch labelled with the conditions  $(\varphi_i)_{i \in \mathbb{N}}$  such that no prefix  $(\varphi_i)_{0 \leq i \leq k}$  ( $k \in \mathbb{N}$ ) is inconsistent. Since Simplex together with Bland’s rule specifies only a finite number of case distinctions for every pivoting step and the number of pivoting steps performed in an acyclic computation is bounded by the number of variables and inequalities in the input system, the branch must describe a cyclic Simplex computation. But this is a contradiction, since the pivots are chosen using Bland’s rule and, therefore, no consistent prefix of the branch can contain a cycle. Thus, no infinite branch with only consistent finite prefixes can exist.

Let us summarize our transformational approach to obtaining algorithms for the generalized polyhedron model:

- Start with an existing algorithm for the current polyhedron model.
- Replace every case distinction in the algorithm’s code by a case distinction in the result data structure (by replacing the code for the case distinction by an application of a suitable decision tree constructor like *GeCond*).
- Change the result type from  $\alpha$  to *Tree*  $\alpha$  in the type signature.

To make the principle clearer, we have omitted some technical details (Größlinger, 2003).

*Example.* An application of our generalized Fourier–Motzkin implementation (which uses our generalized Simplex algorithm to eliminate redundant bounds) to the conjunction of systems (3) and (4) yields a system of projections which correspond directly to the loop bounds of the code in Fig. 5. No case distinctions are needed to represent the result.

### 2.2.2. Using quantifier elimination directly to solve some problems

Some problems in the generalized polyhedron model can be formulated as quantifier elimination problems. For example, testing the feasibility of an inequality system  $\psi$  is simply asking the quantifier elimination tool for a quantifier-free equivalent of  $\exists x_1 \dots \exists x_n(\psi)$ .

Computing the lexicographic minimum of an inequality system is an example in which quantifier elimination *with answer* can be used. Given an inequality system  $\psi$ , the following formula:

$$\mu := \psi \wedge \forall y_1 \dots \forall y_n (\psi[x_1 := y_1, \dots, x_n := y_n] \rightarrow (x_1, \dots, x_n) \leq_{\text{lex}} (y_1, \dots, y_n))$$

describes that  $(x_1, \dots, x_n)$  is the finite lexicographic minimum of  $\psi$ 's solution set. Quantifier elimination with answer can now compute a quantifier-free formula (in the parameters) of the form  $\bigvee \chi_i$  which is equivalent to  $\exists x_1 \dots \exists x_n(\mu)$  and, for the different cases  $\chi_i$ , values for  $(x_1, \dots, x_n)$  are computed such that  $\mu$  holds, i.e., the lexicographic minimum is computed.

In order to compare this approach with the generalization of existing algorithms (Section 2.2.1), we used a generalized Simplex algorithm for computing the lexicographic minimum. It turns out that the generalized Simplex is by orders of magnitude more efficient than the direct encoding described in this section. The reason is that the Simplex algorithm is tuned to its application domain, whereas quantifier elimination cannot make use of domain-specific information.

A second drawback of using quantifier elimination directly is that not all problems in the polyhedron model can be described easily this way. For example, it is not guaranteed that the projection of a polyhedron with non-linear parameters can be described by a conjunction. Consider the following formula:

$$0 \leq x \leq y \quad \wedge \quad p \cdot y \geq 0 \quad \wedge \quad y \leq 10$$

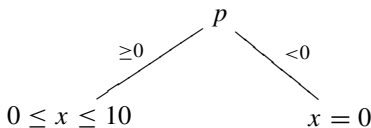
with non-linear parameter  $p$ . If we assume nothing about  $p$ , REDLOG's answer to the question  $\exists y (0 \leq x \leq y \wedge p \cdot y \geq 0 \wedge y \leq 10)$  is:

$$(x = 0 \wedge p \neq 0) \quad \vee \quad (0 \leq x \leq 10 \wedge p \geq 0). \quad (5)$$

This formula cannot be written as a conjunction of atomic formulas. To generate code, the conditions on the loop indices must be joined by conjunctions. The given problem can be solved by applying our generalized Fourier–Motzkin to:

$$0 \leq x \leq y \quad \wedge \quad p \cdot y \geq 0 \quad \wedge \quad y \leq 10$$

to eliminate  $y$ :

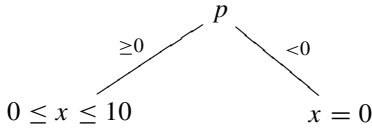


The question now is: can we find the conditions  $p \geq 0$  and  $p < 0$  and the respective solutions  $0 \leq x \leq 10$  and  $x = 0$  from the quantifier elimination's answer to  $\exists y (0 \leq x \leq y \wedge p \cdot y \geq 0 \wedge y \leq 10)$ ?

In this example, one can find a suitable case distinction by extracting the two conditions  $p \neq 0$  and  $p \geq 0$  from the two disjuncts of (5). This yields a total of four possible cases:

1.	$p \neq 0 \wedge p \geq 0$ :	$x = 0 \vee x \geq 0$
2.	$p \neq 0 \wedge \neg(p \geq 0)$ :	$x = 0$
3.	$\neg(p \neq 0) \wedge p \geq 0$ :	$x \geq 0$
4.	$\neg(p \neq 0) \wedge \neg(p \geq 0)$ :	$x \in \emptyset$

After simplification (Case 4 is impossible and Cases 1 and 3 can be combined since both have  $x \geq 0$  as solution) and computing conjunctive equivalents for the formulas describing  $x$  in each case, we get, as desired, the following decision tree:



*Example.* Instead of using our generalized Fourier–Motzkin, one can also use REDLOG and SLFQ to compute the tiled target code (Fig. 5), because we have proved that the projections are equivalent to conjunctions of atomic formulas (Größlinger et al., 2004). As described in Section 2, for every loop index, REDLOG is used to project the input system onto the respective prefix of the loop indices. The resulting formulas are fed to SLFQ to simplify them to a conjunction of atomic formulas from which the loop bounds can be derived. Unfortunately, 2 GB of heap are not enough for SLFQ to simplify the formula describing the innermost loop with the bounds for  $p$ . We can also compute the required bounds using the simpler technique of checking, for every atomic formula in REDLOG's result, whether it is implied by the whole formula. In our example, this works for all loop indices, including  $p$ .

### 3. Conclusion

Automatic, model-based parallelization requires powerful mathematical tools, e.g., for optimization, projection or feasibility tests. The traditional limitation in this research area is the linearity of the expression, but this is too restrictive even for problems which are quite simple but relevant in practice.

The basis for the solution presented here is a powerful and, at the same time, efficiently implemented mathematical tool. Quantifier elimination, e.g., the method by Weispfenning as implemented in the REDLOG package of REDUCE, provides us with the necessary functionality. We discovered that generalizing an existing algorithm (like the Simplex algorithm or Fourier–Motzkin elimination) using quantifier elimination to simplify the resulting decision trees usually performs better than expressing the problem as a logical formula and using quantifier elimination alone to solve it. Combining quantifier elimination with existing algorithms enables a much desired feature, namely non-linear parameters. This way, quantifier elimination is making a significant contribution to bringing parallelization in the polyhedron model closer to the realm of practice.

## Acknowledgements

This work has been funded in part by exchange grants from the German Academic Exchange Service (DAAD) in its Procope programme, and from the Bavarian–French University Center (BFHZ-CCUFB). Our exchange partners were Paul Feautrier and his group; we thank them for discussions. We also thank the reviewers for their useful comments.

## References

- Amiranoff, Pierre, 2004. An automata-theoretic modelization of instancewise program analysis: Transducers as mappings from instances to memory locations. Ph.D. Thesis. Conservatoire National des Arts et Métiers, December 2004.
- Andonov, Rumen, Balev, Stefan, Rajopadhye, Sanjay, Yanev, Nicola, 2001. Optimal semi-oblique tiling. In: Proceedings of the 13th Annual ACM Symposium on Parallel Algorithms and Architectures. SPAA 2001. ACM Press, pp. 153–162. Extended version available as technical report: IRISA, no. 1392, December 2001.
- Bastoul, Cédric, 2004. Code generation in the polyhedral model is easier than you think. In: Proc. 13th Int. Conf. on Parallel Architectures and Compilation Techniques. PACT 2004. IEEE Computer Society Press, pp. 7–16.
- Bland, Robert G., 1977. New finite pivoting rules for the simplex method. *Mathematics of Operations Research* 2, 103–107.
- Brown, Christopher W., 2002. The SLFQ system, <http://www.cs.usna.edu/~qepcad/SLFQ/Home.html>.
- Brown, Christopher W., 2004. QEPCAD quantifier elimination tool. <http://www.cs.usna.edu/~qepcad/B/QEPCAD.html>.
- Chernikova, N.V., 1968. Algorithm for discovering the set of all solutions of a linear programming problem. *U.S.S.R. Computational Mathematics and Mathematical Physics* 8 (6), 282–293.
- Collard, Jean-François, Griebel, Martin, 1999. A precise fixpoint reaching definition analysis for arrays. In: Carter, Larry, Ferrante, Jeanne (Eds.), *Languages and Compilers for Parallel Computing*, 12th International Workshop. LCPC'99. In: LNCS, vol. 1863. Springer-Verlag, pp. 286–302.
- Dion, Michèle, Robert, Yves, 1995. Mapping affine loop nests: New results. In: Hertzberger, Bob, Serazzi, Giuseppe (Eds.), *High-Performance Computing & Networking. HPCN'95*. In: LNCS, vol. 919. Springer-Verlag, pp. 184–189.
- Dolzmann, Andreas, Sturm, Thomas, 1997. REDLOG: Computer algebra meets computer logic. *ACM SIGSAM Bulletin* 31, 2–9.
- Feautrier, Paul, 1991. Dataflow analysis of array and scalar references. *International Journal of Parallel Programming* 20 (1), 23–53.
- Feautrier, Paul, 1992a. Some efficient solutions to the affine scheduling problem. Part I. One-dimensional time. *International Journal of Parallel Processing* 21 (5), 313–348.
- Feautrier, Paul, 1992b. Some efficient solutions to the affine scheduling problem. Part II. Multidimensional time. *International Journal of Parallel Processing* 21 (6), 389–420.
- Feautrier, Paul, 1994. Toward automatic distribution. *Parallel Processing Letters* 4 (3), 233–244.
- Feautrier, Paul, 1996. Automatic parallelization in the polytope model. In: Perrin, Guy-René, Darte, Alain (Eds.), *The Data Parallel Programming Model*. In: LNCS, vol. 1132. Springer-Verlag, pp. 79–103.
- Feautrier, Paul, Collard, Jean-François, Bastoul, Cédric, 2003. Solving systems of affine (in)equalities: PIP's user's guide. Technical Report. PRISM, University of Versailles.
- Griebel, Martin, Faber, Peter, Lengauer, Christian, 2004. Space–time mapping and tiling — a helpful combination. In: Proc. 9th Workshop on Compilers for Parallel Computers. CPC 2001. Concurrency and Computation: Practice and Experience 16 (2–3), 221–246.
- Griebel, Martin, Lengauer, Christian, 1996. The loop parallelizer LooPo. In: Michael, Gerndt (Ed.), *Proc. Sixth Workshop on Compilers for Parallel Computers. CPC'96*. In: *Konferenzen des Forschungszentrums Jülich*, vol. 21. pp. 311–320. Forschungszentrum Jülich.
- Größlinger, Armin, 2003. Extending the Polyhedron Model to Inequality Systems with Non-linear Parameters using quantifier elimination. Diploma Thesis. Universität Passau, <http://www.infosun.fmi.uni-passau.de/cl/arbeiten/groesslinger.ps.gz>.
- Größlinger, Armin, Griebel, Martin, Lengauer, Christian, 2004. Introducing non-linear parameters to the polyhedron model. In: Gerndt Michael, Kereku Edmond (Eds.), *Proc. 11th Workshop on Compilers for Parallel Computers. CPC 2004, Research Report Series, LRR-TUM, Technische Universität München*, pp. 1–12.
- Hearn, Anthony C., 2004. REDUCE computer algebra system, <http://www.reduce-algebra.com>.
- Lengauer, Christian, 1993. Loop parallelization in the polytope model. In: Best, Eike (Ed.), *CONCUR'93*. In: LNCS, vol. 715. Springer-Verlag, pp. 398–416.

- Loos, Rüdiger, Weispfenning, Volker, 1993. Applying linear quantifier elimination. In: Computational Quantifier Elimination. The Computer Journal 36 (5), 450–462 (special issue).
- Pugh, William, 1992. The Omega test: A fast and practical integer programming algorithm for dependence analysis. Communications of the ACM 35 (8), 102–114.
- Pugh, William, Wonnacott, Dave, 1992. Eliminating false data dependences using the Omega test. In: Proc. ACM SIGPLAN Conf. on Programming Language Design and Implementation. PLDI'92. ACM SIGPLAN Notices 27 (7), 140–151.
- Quilleré, Fabien, Rajopadhye, Sanjay, Wilde, Doran, 2000. Generation of efficient nested loops from polyhedra. International Journal of Parallel Processing 28 (5), 469–498.
- Weispfenning, Volker, 1988. The complexity of linear problems in fields. Journal of Symbolic Computation 5 (1–2), 3–27.
- Weispfenning, Volker, 1994a. Parametric linear and quadratic optimization by elimination. Technical Report MIP-9404. Universität Passau.
- Weispfenning, Volker, 1994b. Quantifier elimination for real algebra — the cubic case. In: Proc. Int. Symp. on Symbolic and Algebraic Computation. ISSAC'94. ACM Press, pp. 258–263.
- Weispfenning, Volker, 1997. Quantifier elimination for real algebra — the quadratic case and beyond. Applicable Algebra in Engineering Communication and Computing 8 (2), 85–101.